

Vectorized Implementation for Linear Regression

Quang-Vinh Dinh
Ph.D. in Computer Science

Outline

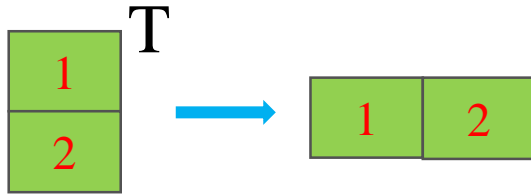
- **Review on One-sample Training**
- **Vectorize the Linear Regression (m-sample)**
- **Vectorize the Linear Regression (N-sample)**
- **Proofs of Some Matrix Properties**

Review

Transpose

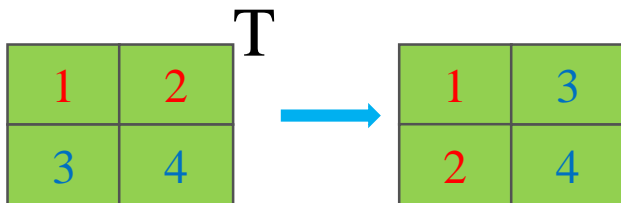
$$\vec{v} = \begin{bmatrix} v_1 \\ \dots \\ v_n \end{bmatrix}$$

$$\vec{v}^T = [v_1 \ \dots \ v_n]$$



$$A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \dots & \dots & \dots \\ a_{m1} & \dots & a_{mn} \end{bmatrix}$$

$$A^T = \begin{bmatrix} a_{11} & \dots & a_{m1} \\ \dots & \dots & \dots \\ a_{1n} & \dots & a_{mn} \end{bmatrix}$$

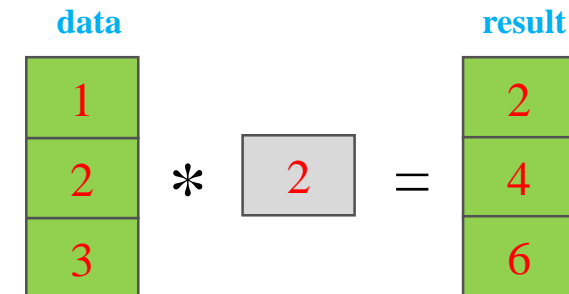


```
2 import numpy as np
3
4 # create data
5 data = np.array([1,2,3])
6 factor = 2
7
8 # broadcasting
9 result_multiplication = data*factor
```

```
[1 2 3]
[2 4 6]
```

Multiply with a number

$$\alpha \vec{u} = \alpha \begin{bmatrix} u_1 \\ \dots \\ u_n \end{bmatrix} = \begin{bmatrix} \alpha u_1 \\ \dots \\ \alpha u_n \end{bmatrix}$$



Review

Feature	Label
area	price
6.7	9.1
4.6	5.9
3.5	4.6
5.5	6.7
x	y

1) Pick a sample (x, y) from training data

2) Compute the output \hat{y}

$$\hat{y} = wx + b$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative

$$\frac{\partial L}{\partial w} = 2x(\hat{y} - y)$$

$$\frac{\partial L}{\partial b} = 2(\hat{y} - y)$$

5) Update parameters

$$w = w - \eta \frac{\partial L}{\partial w}$$

$$b = b - \eta \frac{\partial L}{\partial b}$$

η is learning rate

Traditional

$$\hat{y} = wx + b \quad x = \begin{bmatrix} 1 \\ x \end{bmatrix} \quad \theta = \begin{bmatrix} b \\ w \end{bmatrix}$$

$$\theta = \begin{bmatrix} b \\ w \end{bmatrix} \rightarrow \theta^T = [b \ w]$$

$$\hat{y} = wx + b1 = [b \ w] \begin{bmatrix} 1 \\ x \end{bmatrix} = \theta^T x$$

dot product

Review

1) Pick a sample (x, y) from training data

2) Compute the output \hat{y}

$$\hat{y} = wx + b$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative

$$\frac{\partial L}{\partial w} = 2x(\hat{y} - y)$$

$$\frac{\partial L}{\partial b} = 2(\hat{y} - y)$$

5) Update parameters

$$w = w - \eta \frac{\partial L}{\partial w}$$

$$b = b - \eta \frac{\partial L}{\partial b}$$

η is learning rate

Traditional

$$\hat{y} = wx + b \quad x = \begin{bmatrix} 1 \\ x \end{bmatrix} \quad \theta = \begin{bmatrix} b \\ w \end{bmatrix}$$

$$\hat{y} = \theta^T x$$

$$L(\hat{y}, y) = (\hat{y} - y)^2$$

numbers

What will we do?

Review

Traditional

1) Pick a sample (x, y) from training data

2) Compute the output \hat{y}

$$\hat{y} = wx + b$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative

$$\frac{\partial L}{\partial w} = 2x(\hat{y} - y) \quad \frac{\partial L}{\partial b} = 2(\hat{y} - y)$$

5) Update parameters

$$w = w - \eta \frac{\partial L}{\partial w} \quad b = b - \eta \frac{\partial L}{\partial b}$$

η is learning rate

$$\hat{y} = wx + b \quad x = \begin{bmatrix} 1 \\ x \end{bmatrix} \quad \theta = \begin{bmatrix} b \\ w \end{bmatrix}$$

$$\left\{ \begin{array}{l} \frac{\partial L}{\partial b} = 2(\hat{y} - y) = 2 \times (\hat{y} - y) \times 1 \\ \frac{\partial L}{\partial w} = 2x(\hat{y} - y) = 2 \times (\hat{y} - y) \times x \end{array} \right.$$

$$\begin{bmatrix} 2 \times (\hat{y} - y) \times 1 \\ 2 \times (\hat{y} - y) \times x \end{bmatrix} = \underbrace{2(\hat{y} - y)}_{\text{common factor}} \begin{bmatrix} 1 \\ x \end{bmatrix} = 2(\hat{y} - y)\mathbf{x} = \begin{bmatrix} \frac{\partial L}{\partial b} \\ \frac{\partial L}{\partial w} \end{bmatrix} = L'_{\theta} \quad \rightarrow \quad L'_{\theta} = 2\mathbf{x}(\hat{y} - y)$$

Review

1) Pick a sample (x, y) from training data

2) Compute the output \hat{y}

$$\hat{y} = wx + b$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative

$$\frac{\partial L}{\partial w} = 2x(\hat{y} - y) \quad \frac{\partial L}{\partial b} = 2(\hat{y} - y)$$

5) Update parameters

$$w = w - \eta \frac{\partial L}{\partial w} \quad b = b - \eta \frac{\partial L}{\partial b}$$

Traditional

$$\hat{y} = wx + b$$

$$x = \begin{bmatrix} 1 \\ x \end{bmatrix}$$

$$\theta = \begin{bmatrix} b \\ w \end{bmatrix}$$

$$L'_\theta = \begin{bmatrix} \frac{\partial L}{\partial b} \\ \frac{\partial L}{\partial w} \end{bmatrix}$$

$$\left\{ \begin{array}{l} b = b - \eta \frac{\partial L}{\partial b} \\ w = w - \eta \frac{\partial L}{\partial w} \\ \theta = \theta - \eta L'_\theta \end{array} \right.$$

$$\rightarrow \theta = \theta - \eta L'_\theta$$

Review

1) Pick a sample (x, y) from training data

2) Compute the output \hat{y}

$$\hat{y} = wx + b$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative

$$\frac{\partial L}{\partial w} = 2x(\hat{y} - y)$$

$$\frac{\partial L}{\partial b} = 2(\hat{y} - y)$$

5) Update parameters

$$w = w - \eta \frac{\partial L}{\partial w}$$

$$b = b - \eta \frac{\partial L}{\partial b}$$

η is learning rate

Traditional

1) Pick a sample (x, y) from training data

2) Compute output \hat{y}

$$\hat{y} = \boldsymbol{\theta}^T \mathbf{x} = \mathbf{x}^T \boldsymbol{\theta}$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative

$$L'_{\boldsymbol{\theta}} = 2\mathbf{x}(\hat{y} - y)$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta L'_{\boldsymbol{\theta}}$$

η is learning rate

Vectorized

$$\hat{y} = wx + b$$

$$\mathbf{x} = \begin{bmatrix} 1 \\ x \end{bmatrix} \quad \boldsymbol{\theta} = \begin{bmatrix} b \\ w \end{bmatrix}$$

Feature	Label
area	price
6.7	9.1
4.6	5.9
3.5	4.6
5.5	6.7
x	y

1

$$\mathbf{x} = \begin{bmatrix} 1 \\ x \end{bmatrix} = \begin{bmatrix} 1 \\ 6.7 \end{bmatrix}$$

Given $\boldsymbol{\theta} = \begin{bmatrix} b \\ w \end{bmatrix} = \begin{bmatrix} 0.049 \\ -0.34 \end{bmatrix}$
 $\eta = 0.01$

1) Pick a sample (\mathbf{x}, y) from training data

2) Compute output \hat{y}

$$\hat{y} = \boldsymbol{\theta}^T \mathbf{x} = \mathbf{x}^T \boldsymbol{\theta}$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

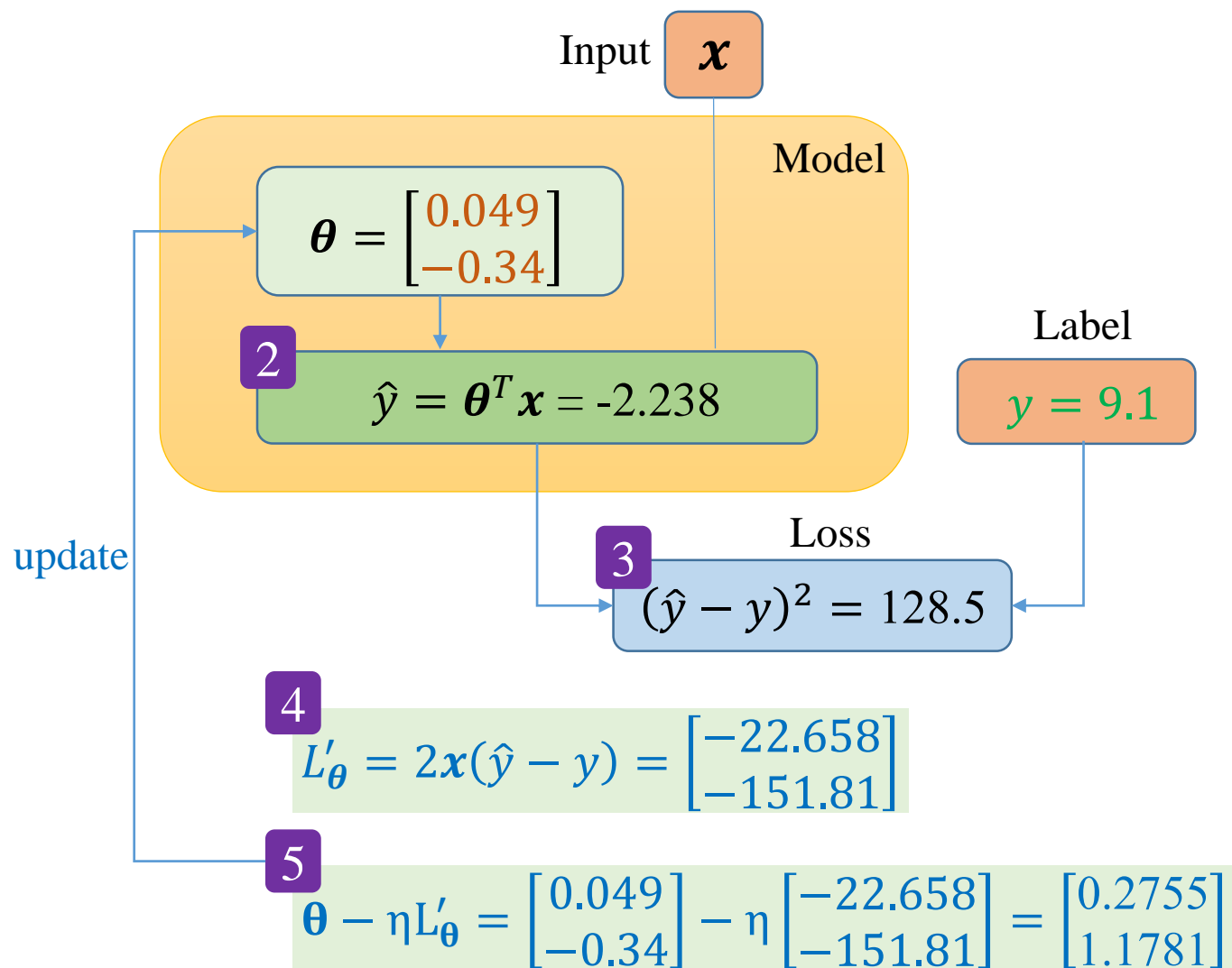
4) Compute derivative

$$L'_{\boldsymbol{\theta}} = 2\mathbf{x}(\hat{y} - y)$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta L'_{\boldsymbol{\theta}}$$

η is learning rate



Objective

❖ Implementation (vectorization using numpy)

1) Pick a sample (x, y) from training data

2) Compute output \hat{y}

$$\hat{y} = \theta^T x = x^T \theta$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative

$$L'_\theta = 2x(\hat{y} - y)$$

5) Update parameters

$$\theta = \theta - \eta L'_\theta$$

η is learning rate

```
1 import numpy as np
2
3 # forward
4 def predict(x, theta):
5     return x.dot(theta)
6
7 # compute gradient
8 def gradient(y_hat, y, x):
9     dtheta = 2*x*(y_hat-y)
10
11     return dtheta
12
13 # update weights
14 def update_weight(theta, lr, dtheta):
15     dtheta_new = theta - lr*dtheta
16
17     return dtheta_new
```

$$\hat{y} = wx + b$$

$$\mathbf{x} = \begin{bmatrix} 1 \\ x \end{bmatrix} \quad \boldsymbol{\theta} = \begin{bmatrix} b \\ w \end{bmatrix}$$

Feature	Label
area	price
6.7	9.1
4.6	5.9
3.5	4.6
5.5	6.7
x	y

1

$$\mathbf{x} = \begin{bmatrix} 1 \\ x \end{bmatrix} = \begin{bmatrix} 1 \\ 3.5 \end{bmatrix}$$

Given $\boldsymbol{\theta} = \begin{bmatrix} b \\ w \end{bmatrix} = \begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}$
 $\eta = 0.01$

1) Pick a sample (x, y) from training data

2) Compute output \hat{y}

$$\hat{y} = \boldsymbol{\theta}^T \mathbf{x} = \mathbf{x}^T \boldsymbol{\theta}$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

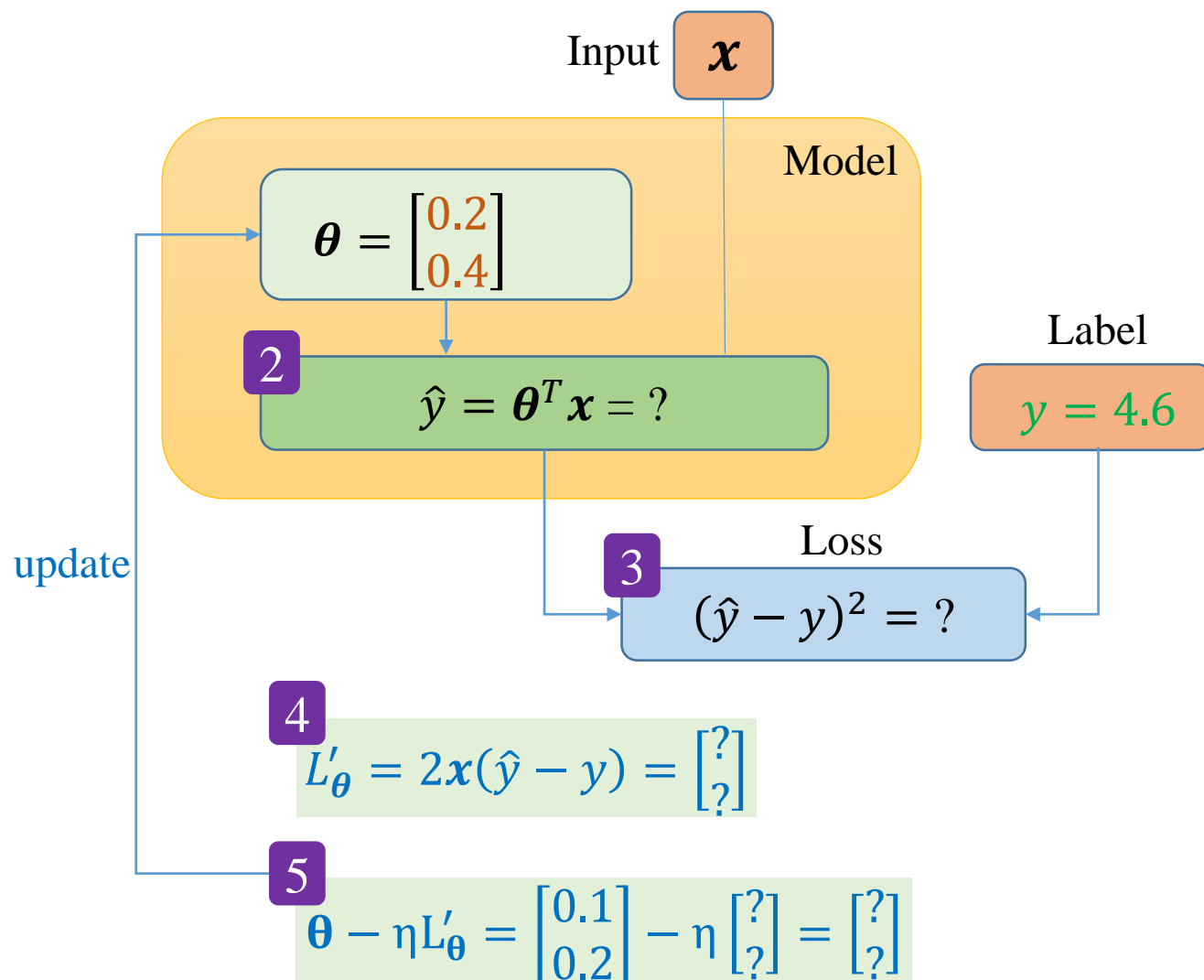
4) Compute derivative

$$L'_{\boldsymbol{\theta}} = 2\mathbf{x}(\hat{y} - y)$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta L'_{\boldsymbol{\theta}}$$

η is learning rate



```

3 data = np.genfromtxt('data.csv', delimiter=',')
4 N = 4
5
6 areas = data[:, 0].reshape(N, 1)
7 prices = data[:, 1].reshape(N,)
8
9 # vector [1, area]
10 features = np.hstack([np.ones((N,1)), areas])

```

```

1 # forward
2 def predict(x, theta):
3     return x.T.dot(theta)
4
5 # compute gradient
6 def gradient(y_hat, y, x):
7     dtheta = 2*x*(y_hat-y)
8     return dtheta
9
10 # update weights
11 def update_weight(theta, lr, dtheta):
12     dtheta_new = theta - lr*dtheta
13     return dtheta_new

```

```

1 lr = 0.01
2 epoch_max = 10
3
4 # [b, w]
5 theta = np.array([0.049, -0.34])
6
7 for epoch in range(epoch_max):
8     for i in range(N):
9         # get a sample
10         x = features[i,:]
11         y = prices[i]
12
13         # predict y_hat
14         y_hat = predict(x, theta)
15
16         # compute loss
17         loss = (y_hat-y)*(y_hat-y)
18
19         # compute gradient
20         dtheta = gradient(y_hat, y, x)
21
22         # update weights
23         theta = update_weight(theta, lr, dtheta)

```

Advertising Problem

Features			Label
TV	↕ Radio	↕ Newspaper	↕ Sales
230.1	37.8	69.2	22.1
44.5	39.3	45.1	10.4
17.2	45.9	69.3	12
151.5	41.3	58.5	16.5
180.8	10.8	58.4	17.9

Advertising data

if TV=55.0, Radio=34.0,
and Newspaper=62.0,
price=?

$$\hat{y} = w_1 * x_1 + w_2 * x_2 + w_3 * x_3 + b$$

$$x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\theta = \begin{bmatrix} b \\ w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

1) Pick a sample (x, y) from training data



2) Compute output \hat{y}



$$\hat{y} = \theta^T x = x^T \theta$$

3) Compute loss



$$L = (\hat{y} - y)^2$$

4) Compute derivative



$$L'_\theta = 2x(\hat{y} - y)$$

5) Update parameters

$$\theta = \theta - \eta L'_\theta$$

η is learning rate

Outline

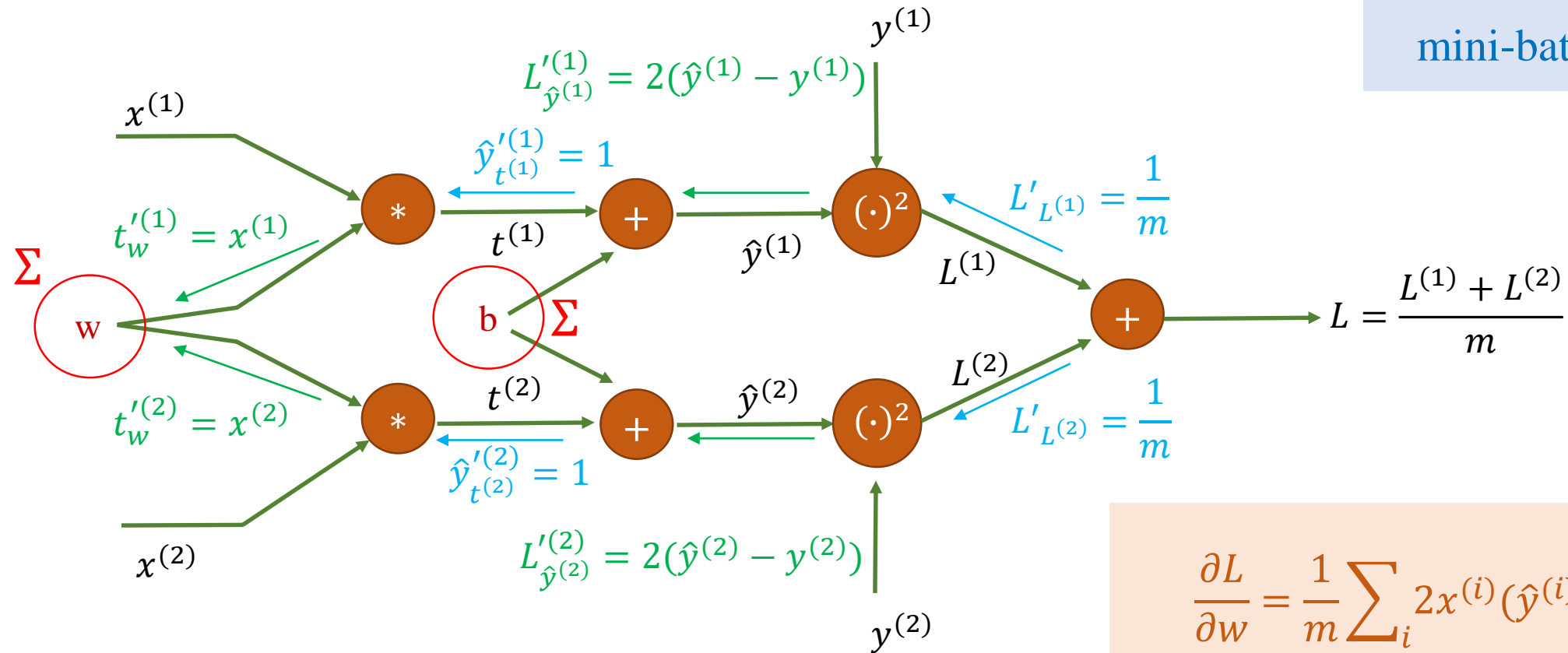
- **Review on One-sample Training**
- **Vectorize the Linear Regression (m-sample)**
- **Vectorize the Linear Regression (N-sample)**
- **Proofs of Some Matrix Properties**

Linear Regression (m-samples)

Feature	Label
area	price
6.7	9.1
4.6	5.9

❖ Compute derivate for w and b

mini-batch $m = 2$



$$\frac{\partial L}{\partial w} = \frac{1}{m} \sum_i 2x^{(i)} (\hat{y}^{(i)} - y^{(i)})$$

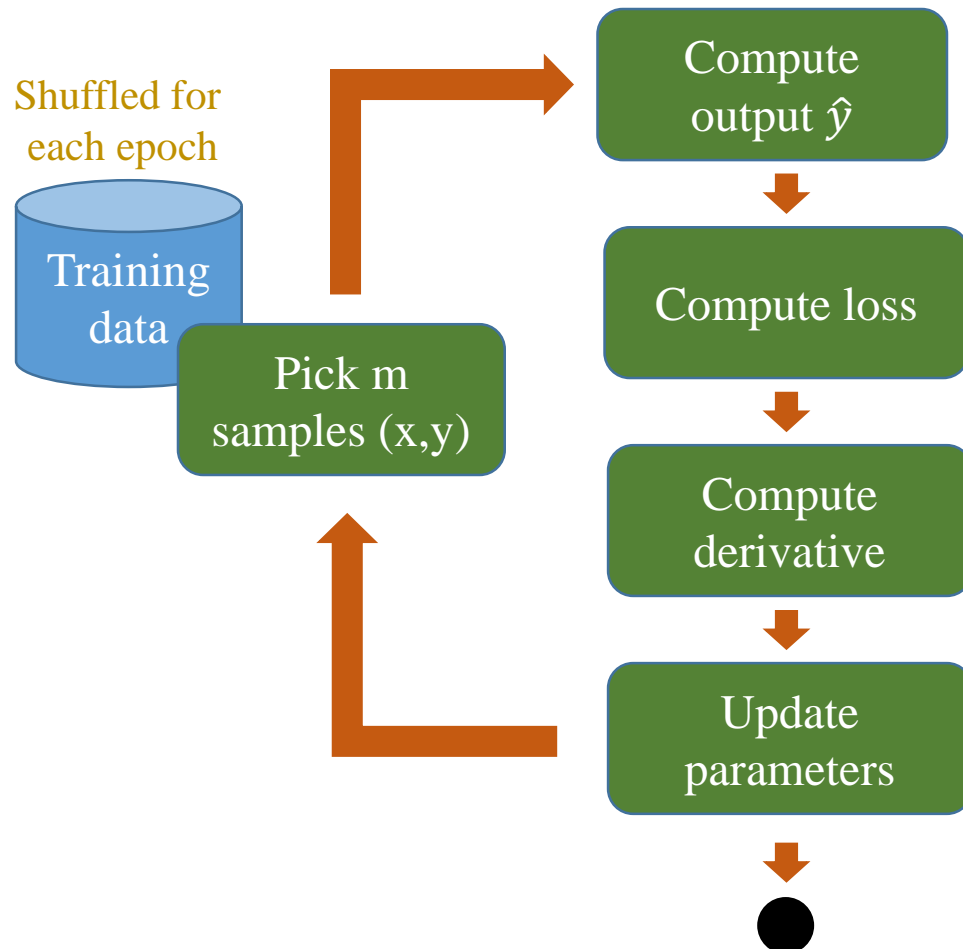
$$\frac{\partial L}{\partial b} = \frac{1}{m} \sum_i 2(\hat{y}^{(i)} - y^{(i)})$$

Linear Regression (m-samples)

	Feature	Label
	area	price
	6.7	9.1
	4.6	5.9

❖ House price prediction

❖ m-sample training ($1 < m < N$)



1) Pick m samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = wx^{(i)} + b \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L^{(i)} = (\hat{y}^{(i)} - y^{(i)})^2 \quad \text{for } 0 \leq i < m$$

4) Compute derivative

$$L'_w{}^{(i)} = 2x^{(i)}(\hat{y}^{(i)} - y^{(i)})$$
$$L'_b{}^{(i)} = 2(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$w = w - \eta \frac{\sum_i L'_w{}^{(i)}}{m}$$

$$b = b - \eta \frac{\sum_i L'_b{}^{(i)}}{m}$$

Learning rate η

Linear Regression (m-samples)

Feature		Label
area		price
6.7		9.1
4.6		5.9

1) Pick a sample (\mathbf{x}, y) from training data

2) Compute output \hat{y}

$$\hat{y} = \boldsymbol{\theta}^T \mathbf{x} = \mathbf{x}^T \boldsymbol{\theta}$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative

$$L'_{\boldsymbol{\theta}} = 2\mathbf{x}(\hat{y} - y)$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta L'_{\boldsymbol{\theta}}$$

η is learning rate

1) Pick m samples $(\mathbf{x}^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} = (\mathbf{x}^{(i)})^T \boldsymbol{\theta} \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L^{(i)} = (\hat{y}^{(i)} - y^{(i)})^2 \quad \text{for } 0 \leq i < m$$

4) Compute derivative

$$L'^{(i)}_{\boldsymbol{\theta}} = 2\mathbf{x}^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{\sum_i L'^{(i)}_{\boldsymbol{\theta}}}{m} \quad \eta \text{ is learning rate}$$

Linear Regression (m-samples)

1) Pick m samples $(\mathbf{x}^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} = (\mathbf{x}^{(i)})^T \boldsymbol{\theta} \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L^{(i)} = (\hat{y}^{(i)} - y^{(i)})^2 \quad \text{for } 0 \leq i < m$$

4) Compute derivative

$$L'_{\boldsymbol{\theta}}^{(i)} = 2\mathbf{x}^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{\sum_i L'_{\boldsymbol{\theta}}^{(i)}}{m} \quad \eta \text{ is learning rate}$$

```
14 for epoch in range(epoch_max):
15     for j in range(0, data_size, m):
16
17         # some variables
18         sum_of_losses = 0
19         gradients = np.zeros((2,))
20         for index in range(j, j+m):
21             # get mini-batch
22             x_i = data[index]
23             y_i = prices[index]
24
25             # predict y_hat_i
26             y_hat_i = x_i.dot(theta)
27
28             # compute loss
29             l_i = (y_hat_i - y_i)*(y_hat_i - y_i)
30
31             # compute gradient
32             gradient_i = x_i*2*(y_hat_i - y_i)
33
34             # accumulate gradients
35             gradients = gradients + gradient_i
36             sum_of_losses = sum_of_losses + l_i
37
38         # normalize
39         sum_of_losses = sum_of_losses/2
40         gradients = gradients/2
```

Linear Regression (m-samples)

❖ Construct formulas

Feature	Label	
area	price	
6.7	9.1	
4.6	5.9	
3.5	4.6	
5.5	6.7	

$$\mathbf{y} = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\mathbf{x} = \begin{bmatrix} 1 & 6.7 \\ 1 & 4.6 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} b \\ w \end{bmatrix} = \begin{bmatrix} 0.049 \\ -0.34 \end{bmatrix}$$

2) Compute output \hat{y}

$$\hat{y} = wx + b$$

$$\mathbf{x} = \begin{bmatrix} 1 & 6.7 \\ 1 & 4.6 \end{bmatrix} = \begin{bmatrix} 1 & x^{(1)} \\ 1 & x^{(2)} \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} b \\ w \end{bmatrix} = \begin{bmatrix} 0.049 \\ -0.34 \end{bmatrix}$$

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \end{bmatrix} = \begin{bmatrix} wx^{(1)} + b \\ wx^{(2)} + b \end{bmatrix} = \begin{bmatrix} 1 & x^{(1)} \\ 1 & x^{(2)} \end{bmatrix} \begin{bmatrix} b \\ w \end{bmatrix} = \mathbf{x}\boldsymbol{\theta} = \begin{bmatrix} -2.229 \\ -1.515 \end{bmatrix}$$

Linear Regression (m-samples)

❖ Construct formulas

Feature	Label	
area	price	
6.7	9.1	
4.6	5.9	
3.5	4.6	
5.5	6.7	

$$\mathbf{y} = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\mathbf{x} = \begin{bmatrix} 1 & 6.7 \\ 1 & 4.6 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} b \\ w \end{bmatrix} = \begin{bmatrix} 0.049 \\ -0.34 \end{bmatrix}$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \end{bmatrix} = \mathbf{x}\boldsymbol{\theta}$$

$$\begin{aligned} L(\hat{\mathbf{y}}, \mathbf{y}) &= \begin{bmatrix} L(\hat{y}^{(1)}, y^{(1)}) \\ L(\hat{y}^{(2)}, y^{(2)}) \end{bmatrix} = \begin{bmatrix} (\hat{y}^{(1)} - y^{(1)})^2 \\ (\hat{y}^{(2)} - y^{(2)})^2 \end{bmatrix} = \begin{bmatrix} (\hat{y}^{(1)} - y^{(1)}) * (\hat{y}^{(1)} - y^{(1)}) \\ (\hat{y}^{(2)} - y^{(2)}) * (\hat{y}^{(2)} - y^{(2)}) \end{bmatrix} \\ &= \begin{bmatrix} (\hat{y}^{(1)} - y^{(1)}) \\ (\hat{y}^{(2)} - y^{(2)}) \end{bmatrix} \odot \begin{bmatrix} (\hat{y}^{(1)} - y^{(1)}) \\ (\hat{y}^{(2)} - y^{(2)}) \end{bmatrix} = (\hat{\mathbf{y}} - \mathbf{y}) \odot (\hat{\mathbf{y}} - \mathbf{y}) \end{aligned}$$

Hadamard product

Linear Regression (m-samples)

❖ Construct formulas

Feature	Label
area	price
6.7	9.1
4.6	5.9
3.5	4.6
5.5	6.7

$$\mathbf{y} = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\mathbf{x} = \begin{bmatrix} 1 & 6.7 \\ 1 & 4.6 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} b \\ w \end{bmatrix} = \begin{bmatrix} 0.049 \\ -0.34 \end{bmatrix}$$

4) Compute derivative

$$\frac{\partial L}{\partial b} = 2(\hat{\mathbf{y}} - \mathbf{y})$$

$$\frac{\partial L}{\partial \mathbf{w}} = 2\mathbf{x}(\hat{\mathbf{y}} - \mathbf{y})$$

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \end{bmatrix} = \mathbf{x}\boldsymbol{\theta}$$

one sample

$$L'_{\boldsymbol{\theta}} = \begin{bmatrix} \frac{\partial L}{\partial b} \\ \frac{\partial L}{\partial \mathbf{w}} \end{bmatrix} = 2(\hat{y} - y) \begin{bmatrix} 1 \\ x \end{bmatrix} = 2(\hat{y} - y)\mathbf{x}$$

two samples

$$2(\hat{\mathbf{y}} - \mathbf{y}) = \begin{bmatrix} 2(\hat{y}^{(1)} - y^{(1)}) \\ 2(\hat{y}^{(2)} - y^{(2)}) \end{bmatrix}$$

$$\mathbf{x} = \begin{bmatrix} 1 & 6.7 \\ 1 & 4.6 \end{bmatrix}$$

Linear Regression (m-samples)

❖ Construct formulas

	Feature	Label	
	area	price	
	6.7	9.1	
	4.6	5.9	
	3.5	4.6	
	5.5	6.7	

$$\mathbf{y} = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\mathbf{x} = \begin{bmatrix} 1 & 6.7 \\ 1 & 4.6 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} b \\ w \end{bmatrix} = \begin{bmatrix} 0.049 \\ -0.34 \end{bmatrix}$$

one sample

$$L'_{\boldsymbol{\theta}} = \begin{bmatrix} \frac{\partial L}{\partial b} \\ \frac{\partial L}{\partial w} \end{bmatrix} = 2(\hat{y} - y) \begin{bmatrix} 1 \\ x \end{bmatrix} = 2(\hat{y} - y)\mathbf{x}$$

propagate and accumulate

two samples

$$\mathbf{k} = 2(\hat{\mathbf{y}} - \mathbf{y}) = \begin{bmatrix} 2(\hat{y}^{(1)} - y^{(1)}) \\ 2(\hat{y}^{(2)} - y^{(2)}) \end{bmatrix}$$

$$\mathbf{x} = \begin{bmatrix} 1 & 6.7 \\ 1 & 4.6 \end{bmatrix}$$

how?

Linear Regression (m-samples)

❖ Construct formulas

	Feature	Label	
	area	price	
	6.7	9.1	
	4.6	5.9	
	3.5	4.6	
	5.5	6.7	

$$\mathbf{y} = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\mathbf{x} = \begin{bmatrix} 1 & 6.7 \\ 1 & 4.6 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} b \\ w \end{bmatrix} = \begin{bmatrix} 0.049 \\ -0.34 \end{bmatrix}$$

two samples

$$\mathbf{k} = 2(\hat{\mathbf{y}} - \mathbf{y}) = \begin{bmatrix} 2(\hat{y}^{(1)} - y^{(1)}) \\ 2(\hat{y}^{(2)} - y^{(2)}) \end{bmatrix}$$

propagate and accumulate

$$\mathbf{x} = \begin{bmatrix} 1 & 6.7 \\ 1 & 4.6 \end{bmatrix}$$

way 1

$$\begin{bmatrix} 2(\hat{y}^{(1)} - y^{(1)}) \\ 2(\hat{y}^{(2)} - y^{(2)}) \end{bmatrix} \begin{bmatrix} 1 & 6.7 \\ 1 & 4.6 \end{bmatrix} = \begin{bmatrix} \frac{\partial L}{\partial b} & \frac{\partial L}{\partial w} \end{bmatrix} = \mathbf{k}^T \mathbf{x}$$

way 2

$$\begin{bmatrix} 1 & 1 \\ 6.7 & 4.6 \end{bmatrix} \begin{bmatrix} 2(\hat{y}^{(1)} - y^{(1)}) \\ 2(\hat{y}^{(2)} - y^{(2)}) \end{bmatrix} = \begin{bmatrix} \frac{\partial L}{\partial b} \\ \frac{\partial L}{\partial w} \end{bmatrix} = \mathbf{x}^T \mathbf{k}$$

Linear Regression (m-samples)

❖ Construct formulas

	Feature	Label	
	area	price	
	6.7	9.1	
	4.6	5.9	
	3.5	4.6	
	5.5	6.7	

$$y = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$x = \begin{bmatrix} 1 & 6.7 \\ 1 & 4.6 \end{bmatrix}$$

$$\theta = \begin{bmatrix} b \\ w \end{bmatrix} = \begin{bmatrix} 0.049 \\ -0.34 \end{bmatrix}$$

5) Update parameters

$$w = w - \eta \frac{\partial L}{\partial w}$$

$$b = b - \eta \frac{\partial L}{\partial b}$$

$$\left\{ \begin{array}{l} b \\ w \\ \theta \end{array} \right\} = \left\{ \begin{array}{l} b \\ w \\ \theta \end{array} \right\} - \eta \left\{ \begin{array}{l} \frac{\partial L}{\partial b} \\ \frac{\partial L}{\partial w} \\ L'_\theta \end{array} \right\}$$

$$L'_\theta = \begin{bmatrix} \frac{\partial L}{\partial b} \\ \frac{\partial L}{\partial w} \end{bmatrix} = x^T k$$

$$L'_\theta = \begin{bmatrix} \frac{\partial L}{\partial b} & \frac{\partial L}{\partial w} \end{bmatrix} = k^T x$$

which one?

$$\rightarrow \theta = \theta - \eta L'_\theta$$

Linear Regression (m-samples)

❖ Example

$$\eta = 0.01$$

$$m = 2$$

Feature	Label
area	price
6.7	9.1
4.6	5.9
3.5	4.6
5.5	6.7

$$\mathbf{x} = \begin{bmatrix} 1 & 6.7 \\ 1 & 4.6 \end{bmatrix}$$

$$\mathbf{y} = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} b \\ w \end{bmatrix} = \begin{bmatrix} 0.049 \\ -0.34 \end{bmatrix}$$

$$\hat{\mathbf{y}} = \mathbf{x}\boldsymbol{\theta} = \begin{bmatrix} -2.229 \\ -1.515 \end{bmatrix}$$

$$\mathbf{L} = (\hat{\mathbf{y}} - \mathbf{y}) \odot (\hat{\mathbf{y}} - \mathbf{y}) = \begin{bmatrix} 128.3 \\ 54.98 \end{bmatrix}$$

$$\mathbf{k} = 2(\hat{\mathbf{y}} - \mathbf{y}) = \begin{bmatrix} -22.658 \\ -14.830 \end{bmatrix}$$

$$\mathbf{L}'_{\boldsymbol{\theta}} = \mathbf{x}^T \mathbf{k} = \begin{bmatrix} 1 \\ 6.7 \end{bmatrix} \begin{bmatrix} 1 \\ 4.6 \end{bmatrix} \begin{bmatrix} -22.658 \\ -14.830 \end{bmatrix} = \begin{bmatrix} -37.488 \\ -220.02 \end{bmatrix}$$

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{\mathbf{L}'_{\boldsymbol{\theta}}}{m}$$

$$= \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix} - \frac{0.01}{2} \begin{bmatrix} -37.488 \\ -220.02 \end{bmatrix} = \begin{bmatrix} 0.236 \\ 0.760 \end{bmatrix}$$

Linear Regression (m-samples)

❖ Formulas

Feature		Label	
	area		price
	6.7		9.1
	4.6		5.9
	3.5		4.6
	5.5		6.7

$$\mathbf{x} = \begin{bmatrix} 1 & 6.7 \\ 1 & 4.6 \end{bmatrix}$$

$$\mathbf{y} = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} b \\ w \end{bmatrix} = \begin{bmatrix} 0.049 \\ -0.34 \end{bmatrix}$$

1) Pick m samples (\mathbf{x}, \mathbf{y}) from training data

2) Compute output $\hat{\mathbf{y}}$

$$\hat{\mathbf{y}} = \mathbf{x}\boldsymbol{\theta}$$

3) Compute loss

$$L(\hat{\mathbf{y}}, \mathbf{y}) = (\hat{\mathbf{y}} - \mathbf{y}) \odot (\hat{\mathbf{y}} - \mathbf{y})$$

4) Compute derivative

$$\mathbf{k} = 2(\hat{\mathbf{y}} - \mathbf{y})$$

$$L'_{\boldsymbol{\theta}} = \mathbf{x}^T \mathbf{k}$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{L'_{\boldsymbol{\theta}}}{m}$$

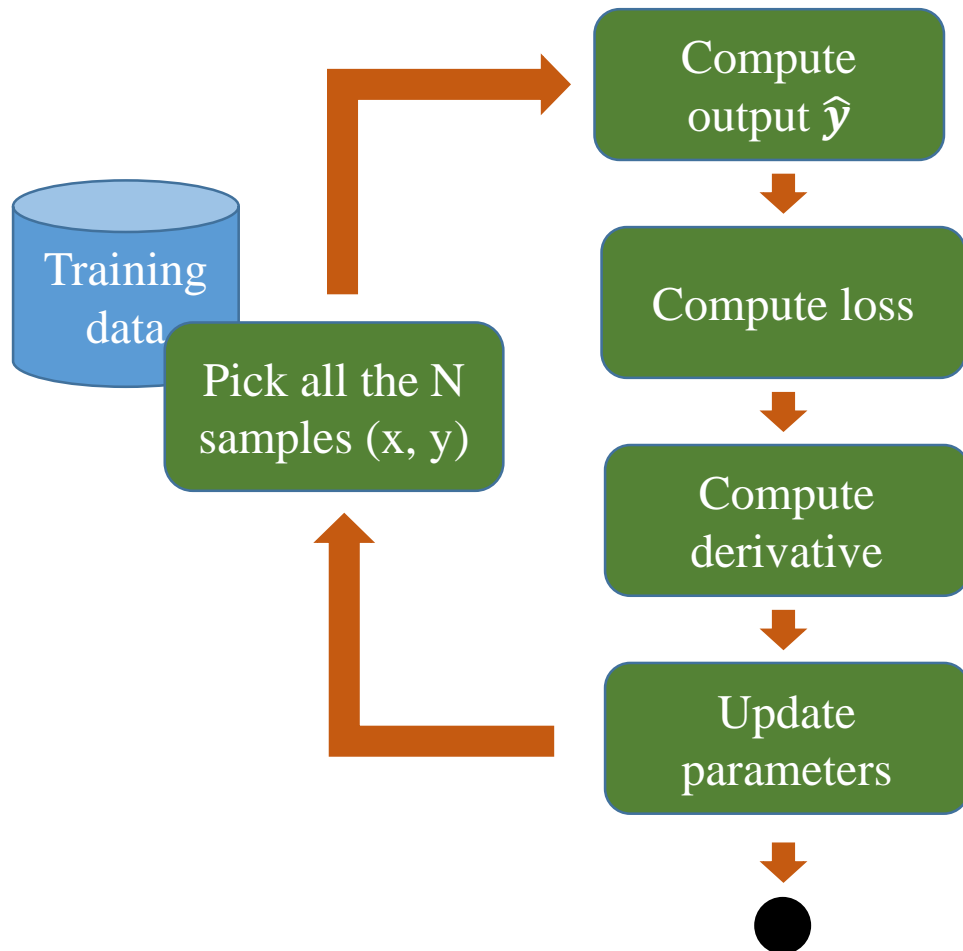
Outline

- **Review on One-sample Training**
- **Vectorize the Linear Regression (m-sample)**
- **Vectorize the Linear Regression (N-sample)**
- **Proofs of Some Matrix Properties**

Linear Regression (m-samples)

❖ House price prediction

❖ N-sample training



1) Pick all the N samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = wx^{(i)} + b \quad \text{for } 0 \leq i < N$$

3) Compute loss

$$L^{(i)} = (\hat{y}^{(i)} - y^{(i)})^2 \quad \text{for } 0 \leq i < N$$

4) Compute derivative

$$L'_w{}^{(i)} = 2x^{(i)}(\hat{y}^{(i)} - y^{(i)})$$
$$L'_b{}^{(i)} = 2(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < N$$

5) Update parameters

$$w = w - \eta \frac{\sum_i L'_w{}^{(i)}}{N}$$

$$b = b - \eta \frac{\sum_i L'_b{}^{(i)}}{N} \quad \text{Learning rate } \eta$$

Linear Regression

Feature	Label
area	price
6.7	9.1
4.6	5.9
3.5	4.6
5.5	6.7

House Price Data

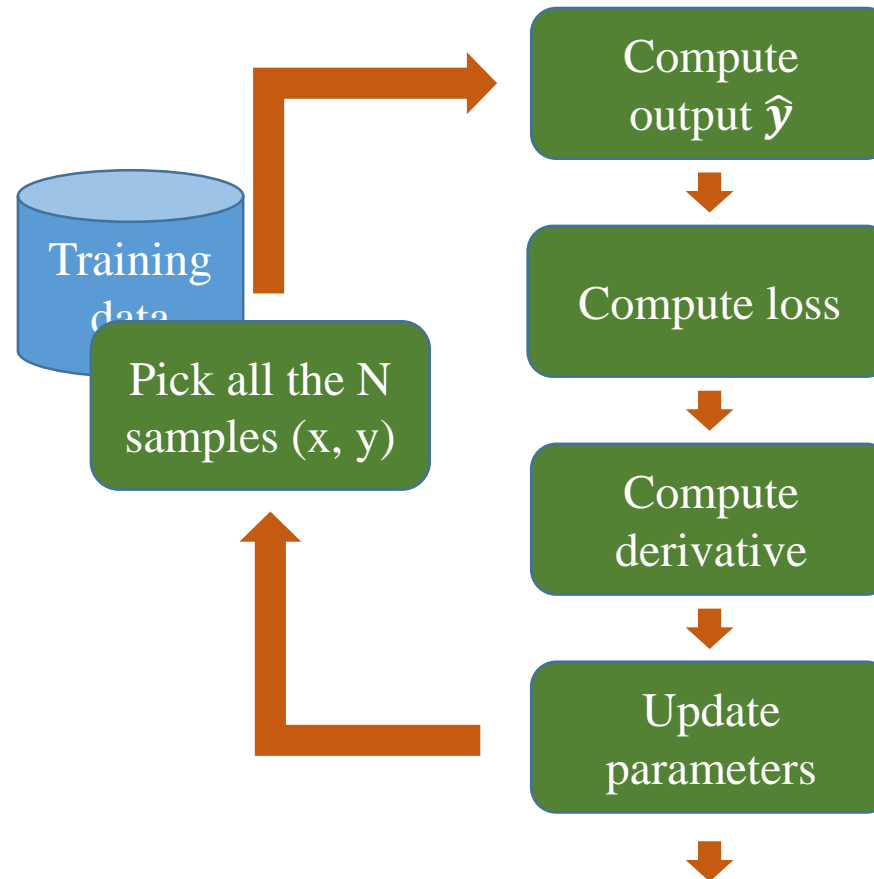
$$x = \begin{bmatrix} 1 & 6.7 \\ 1 & 4.6 \\ 1 & 3.5 \\ 1 & 5.5 \end{bmatrix} \quad y = \begin{bmatrix} 9.1 \\ 5.9 \\ 4.6 \\ 6.7 \end{bmatrix}$$

Model

$$\text{price} = w * \text{area} + b$$
$$\hat{y} = wx + b$$

Parameter Initialization

$$\theta = \begin{bmatrix} b \\ w \end{bmatrix} = \begin{bmatrix} 0.049 \\ -0.34 \end{bmatrix}$$



1) Pick all the N samples from training data

2) Compute output \hat{y}

$$\hat{y} = x\theta$$

3) Compute loss

$$L(\hat{y}, y) = (\hat{y} - y) \odot (\hat{y} - y)$$

4) Compute derivative

$$k = 2(\hat{y} - y)$$

$$L'_{\theta} = x^T k$$

5) Update parameters

$$\theta = \theta - \eta \frac{L'_{\theta}}{N}$$

η is learning rate

Vectorization Approach

Linear Regression

Feature	Label
area	price
6.7	9.1
4.6	5.9
3.5	4.6
5.5	6.7

House Price Data

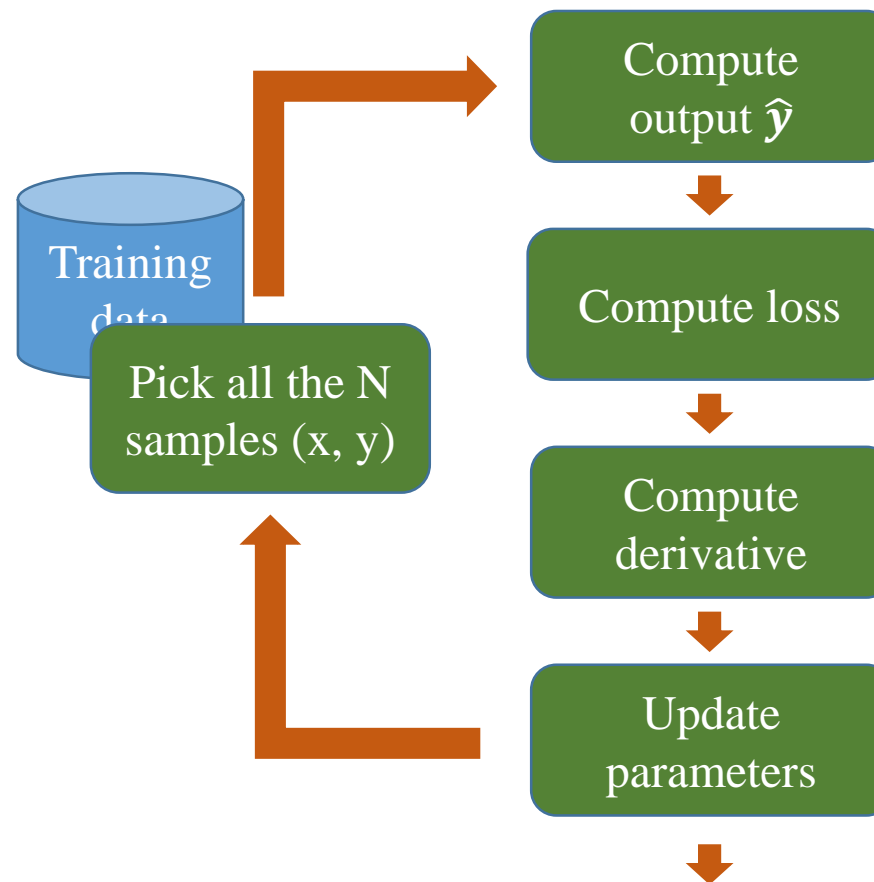
$$x = \begin{bmatrix} 1 & 6.7 \\ 1 & 4.6 \\ 1 & 3.5 \\ 1 & 5.5 \end{bmatrix} \quad y = \begin{bmatrix} 9.1 \\ 5.9 \\ 4.6 \\ 6.7 \end{bmatrix}$$

Model

$$\text{price} = w * \text{area} + b$$
$$\hat{y} = wx + b$$

Parameter Initialization

$$\theta = \begin{bmatrix} b \\ w \end{bmatrix} = \begin{bmatrix} 0.049 \\ -0.34 \end{bmatrix}$$



- 1) Pick all the N samples from training data
- 2) Compute output \hat{y}

$$\hat{y} = x\theta$$

- 3) Compute loss

$$L(\hat{y}, y) = (\hat{y} - y) \odot (\hat{y} - y)$$

- 4) Compute derivative

$$k = 2(\hat{y} - y)$$

$$L'_{\theta} = x^T k$$

- 5) Update parameters

$$\theta = \theta - \eta \frac{L'_{\theta}}{N} \quad \eta \text{ is learning rate}$$

```
for epoch in range(epoch_max):
    y_hat = x.dot(theta)

    loss = np.multiply((y_hat-y), (y_hat-y))
    losses.append(np.mean(loss))

    k = 2*(y_hat-y)
    gradients = x.T.dot(k) / N

    theta = theta - lr*gradients
```

Linear Regression

	Advantages	Disadvantages
1 sample	<ul style="list-style-type: none">Simple to understand and implementFaster learning on some problemsNoisy update is beneficial sometime	<ul style="list-style-type: none">Computationally expensiveNoisy gradient signalConvergence problem
m sample	A balance between the robustness of 1-sample and the efficiency of N-sample	
N sample	<ul style="list-style-type: none">Computationally efficientMore stable error gradientparallel processing	<ul style="list-style-type: none">Premature convergenceMemory problemTraining speed is slower

Outline

- **Review on One-sample Training**
- **Vectorize the Linear Regression (m-sample)**
- **Vectorize the Linear Regression (N-sample)**
- **Proofs of Some Matrix Properties**

$$A \in \mathcal{R}^{m \times n}$$

$$A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \dots & \dots & \dots \\ a_{m1} & \dots & a_{mn} \end{bmatrix}$$

$$B \in \mathcal{R}^{m \times n}$$

$$B = \begin{bmatrix} b_{11} & \dots & b_{1n} \\ \dots & \dots & \dots \\ b_{m1} & \dots & b_{mn} \end{bmatrix}$$

Prove $A + B = B + A$

$$\begin{aligned} A + B &= \begin{bmatrix} (a_{11} + b_{11}) & \dots & (a_{1n} + b_{1n}) \\ \dots & \dots & \dots \\ (a_{m1} + b_{m1}) & \dots & (a_{mn} + b_{mn}) \end{bmatrix} \\ &= \begin{bmatrix} (b_{11} + a_{11}) & \dots & (b_{1n} + a_{1n}) \\ \dots & \dots & \dots \\ (b_{m1} + a_{m1}) & \dots & (b_{mn} + a_{mn}) \end{bmatrix} = B + A \end{aligned}$$

$$A \in \mathcal{R}^{m \times n}$$

$$A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \dots & \dots & \dots \\ a_{m1} & \dots & a_{mn} \end{bmatrix}$$

Prove $(A^T)^T = A$

$$\begin{aligned} (A^T)^T &= \left(\begin{bmatrix} a_{11} & \dots & a_{1n} \\ \dots & \dots & \dots \\ a_{m1} & \dots & a_{mn} \end{bmatrix}^T \right)^T \\ &= \left(\begin{bmatrix} a_{11} & \dots & a_{m1} \\ \dots & \dots & \dots \\ a_{1n} & \dots & a_{mn} \end{bmatrix} \right)^T \\ &= \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \dots & \dots & \dots \\ a_{m1} & \dots & a_{mn} \end{bmatrix} = A \end{aligned}$$

$$A \in \mathcal{R}^{m \times n}$$

$$I \in \mathcal{R}^{n \times n}$$

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}, \quad I = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}$$

Prove $IA = A = AI$

$$AI = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix} \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}$$

$$= \begin{bmatrix} (a_{11}\mathbf{1} + a_{12}0 + \cdots + a_{1n}0) & (a_{11}0 + a_{12}\mathbf{1} + \cdots + a_{1n}0) & \cdots & (a_{11}0 + a_{12}0 + \cdots + a_{1n}\mathbf{1}) \\ (a_{21}\mathbf{1} + a_{22}0 + \cdots + a_{2n}0) & (a_{21}0 + a_{22}\mathbf{1} + \cdots + a_{2n}0) & \cdots & (a_{21}0 + a_{22}0 + \cdots + a_{2n}\mathbf{1}) \\ \vdots & \vdots & \ddots & \vdots \\ (a_{m1}\mathbf{1} + a_{m2}0 + \cdots + a_{mn}0) & (a_{m1}0 + a_{m2}\mathbf{1} + \cdots + a_{mn}0) & \cdots & (a_{m1}0 + a_{m2}0 + \cdots + a_{mn}\mathbf{1}) \end{bmatrix}$$

$$= \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}$$

